

# COMPLEXITY OF SUBREGULAR STRINGSETS

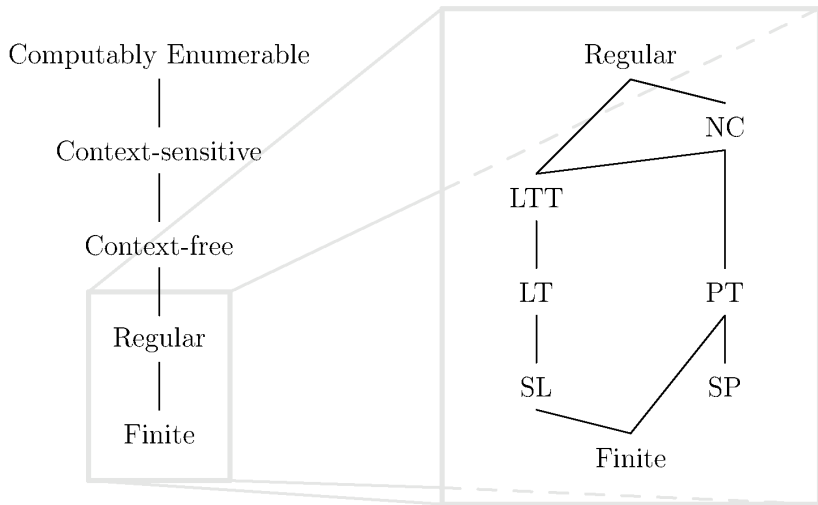
Jeffrey Heinz



Stony Brook University

March 2, 2021

# COMPUTATIONAL COMPLEXITY



# MY OWN INTERESTS IN SUBREGULAR COMPLEXITY...

...began with these observations:

- 1 Many phonotactic patterns are regular.
- 2 Many regular patterns are *not* possible phonotactic patterns.

So *which* regular languages constitute possible phonotactic patterns?

Along the way I learned:

- 1 There is a rich field of study in computer science on *subregular* complexity.
- 2 There is a rich field of study in computer science on *machine learning* of regular languages and transductions.
- 3 There are many applications in linguistics, artificial intelligence, robotic planning and control, model checking,  
...

# MEASURING COMPLEXITY

One idea is that a regular language  $L_1$  is *more complex* than  $L_2$  if

- 1 The smallest DFA recognizing  $L_1$  is larger than the smallest DFA recognizing  $L_2$ .
- 2 The smallest regular expression recognizing  $L_1$  is larger than the smallest regular expression recognizing  $L_2$ .

In contrast to these *intensional* measures

- 1 The subregular classes we study today will provide complexity measures *independent* of the size of such representations.

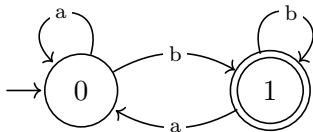
# INTUITIONS?

Consider the four sets of strings below. Any intuitions about which set of strings is more or less complex than any others?

- Strings must end with **b**.
- Strings have odd many **bs**.
- Strings have at least one **b**.
- Strings have exactly one **b**.

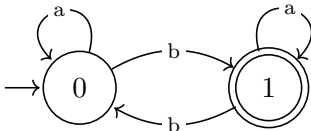
# COMPARING “FINAL-B” WITH “ODD-B”

Strings must end with b.



$(a + b)^*b$

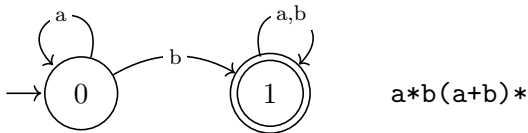
Strings have odd many bs.



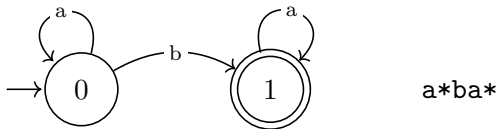
$a^*ba^*(a^*ba^*ba^*)^*$

# COMPARING “AT LEAST ONE B” WITH “EXACTLY ONE B”

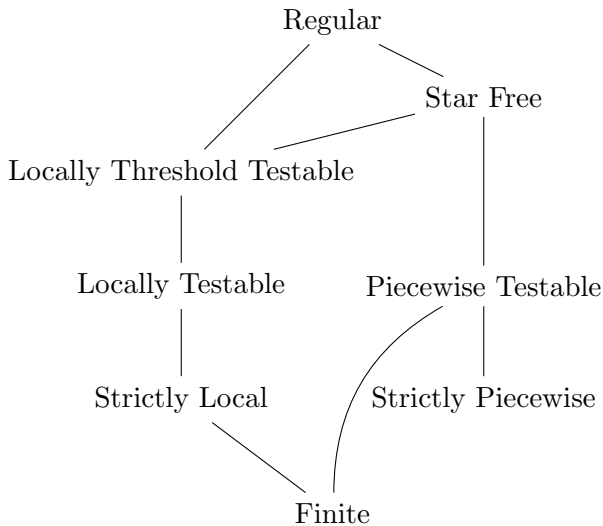
Strings have at least one b.



Strings have exactly one b.

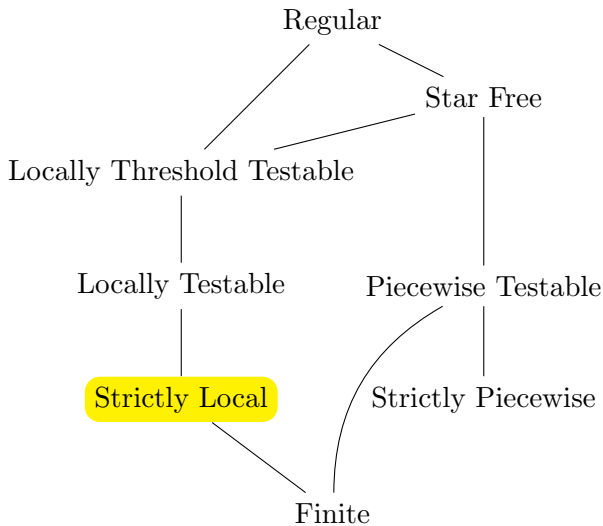


# SUBREGULAR COMPLEXITY





# SUBREGULAR COMPLEXITY



# FACTORS

The notion of *substring*, also called *factor*, is key to understanding the “local” subregular classes.

## k-Factors

- A string  $u$  is a  $k$ -factor (substring of length  $k$ ) of a string  $v$  if there exist strings  $x, y$  such that  $v = xuy$  and  $|u| = k$ . We write  $u \sqsubseteq_k v$ .
- $\mathbf{factor}_k(w)$  equals the set  $\{u \mid u \sqsubseteq_k w\}$  if  $|w| > k$  and equals  $\{w\}$  if  $|w| \leq k$ .
- $\mathbf{factor}_k(S)$  equals the set  $\bigcup_{w \in S} \mathbf{factor}_k(w)$ .

# EXERCISE

Identify the following sets.

- 1  $\text{factor}_2(abcd)$
- 2  $\text{factor}_3(abcd)$
- 3  $\text{factor}_6(abcd)$
- 4  $\text{factor}_2(aaaa)$
- 5  $\text{factor}_2(aaaa, bbbb, caabc)$

# STRICTLY LOCAL LANGUAGES

Let  $\bowtie$  and  $\bowtie$  denote left and right word boundaries.

A language  $L$  is Strictly  $k$ -Local if there is a finite set of strings  $S \subseteq \mathbf{factor}_k(\{\bowtie\} \cdot \Sigma^* \cdot \{\bowtie\})$  such that

$$L = \{w \in \Sigma^* \mid \mathbf{factor}_k(\bowtie w \bowtie) \subseteq S\}$$

A language  $L$  is Strictly Local if there is some  $k$  for which it is Strictly  $k$ -Local. We sometimes write  $\llbracket S \rrbracket$  to denote the extension of  $S$ .

# STRICTLY LOCAL LANGUAGES

## Notes

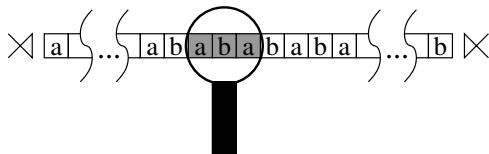
- Elements in  $S$  are called **permissible factors**.
- Elements in  $\mathbf{factor}_k(\{\bowtie\} \cdot \Sigma^* \cdot \{\bowtie\})$  not in  $S$  are called **forbidden factors**.

## Exercises

- 1 Let  $S = \{\bowtie a, ab, ba, a\bowtie\}$ . What is  $\llbracket S \rrbracket$ ?
- 2 For what  $k$  is  $\llbracket S \rrbracket$  Strictly  $k$ -Local?
- 3 Which factors are forbidden?

# STRICTLY LOCAL LANGUAGES: SCANNERS

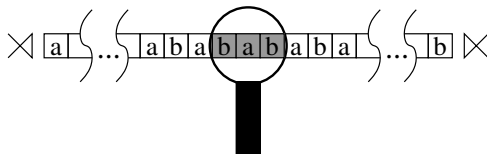
Intuitively, if  $L$  is Strictly  $k$ -Local, then deciding whether a string  $w$  belongs to  $L$  simply requires scanning  $w$  for forbidden substrings. If any is found  $w$  is rejected. If every substring in  $w$  is permissible then  $w$  is accepted.



(McNaughton and Papert 1971, Rogers and Pullum 2011)

# STRICTLY LOCAL LANGUAGES: SCANNERS

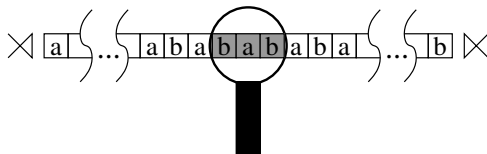
Intuitively, if  $L$  is Strictly  $k$ -Local, then deciding whether a string  $w$  belongs to  $L$  simply requires scanning  $w$  for forbidden substrings. If any is found  $w$  is rejected. If every substring in  $w$  is permissible then  $w$  is accepted.



(McNaughton and Papert 1971, Rogers and Pullum 2011)

# STRICTLY LOCAL LANGUAGES: SCANNERS

Intuitively, if  $L$  is Strictly  $k$ -Local, then deciding whether a string  $w$  belongs to  $L$  simply requires scanning  $w$  for forbidden substrings. If any is found  $w$  is rejected. If every substring in  $w$  is permissible then  $w$  is accepted.



(McNaughton and Papert 1971, Rogers and Pullum 2011)



# MORE EXERCISES

Consider the languages below.

- (1) Strings end with a **b**.
- (2) The second to last symbol in all strings is **b**.
- (3) The third to last symbol in all strings is **b**.

For each, one explain why it is Strictly Local. Assume  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .

# THEOREMS

- 1  $SL_1 \subsetneq SL_2 \dots SL_k \subsetneq SL_{k+1} \dots \subseteq SL$
- 2  $FIN \subsetneq SL$
- 3 For each  $k$ ,  $FIN \not\subseteq SL_k$
- 4 For each  $k$ ,  $SL_k$  is closed under intersection, but neither complement nor union.

# GRAMMAR-INDEPENDENT CHARACTERIZATION OF SL

## Suffix Substitution Closure

A language  $L$  is Strictly Local iff there is a  $k$  such that for all strings  $u_1, v_1, u_2, v_2 \in \Sigma^*$  and for all strings  $x$  of length  $k - 1$  whenever  $u_1xv_1, u_2xv_2 \in L$  then  $u_1xv_2 \in L$ .

In a picture

$$\begin{array}{cccc} & \overbrace{\hspace{1.5cm}}^{k-1} & & \\ u_1 & x & v_1 & \in L \\ u_2 & x & v_2 & \in L \\ \hline u_1 & x & v_2 & \in L \end{array}$$

# PROVING SOME LANGUAGES ARE NOT SL

The SSC helps us prove languages are not Strictly  $k$ -Local and even not Strictly Local for any  $k$ .

- To show  $L$  is not  $SL_k$ , find  $u_1, v_1, u_2, v_2$  and  $x$  of length  $k - 1$  such that

In a picture

$$\begin{array}{cccc} & \underbrace{\hspace{1.5cm}}_{k-1} & & \\ u_1 & x & v_1 & \in L \\ u_2 & x & v_2 & \in L \\ \hline u_1 & x & v_2 & \notin L \end{array}$$

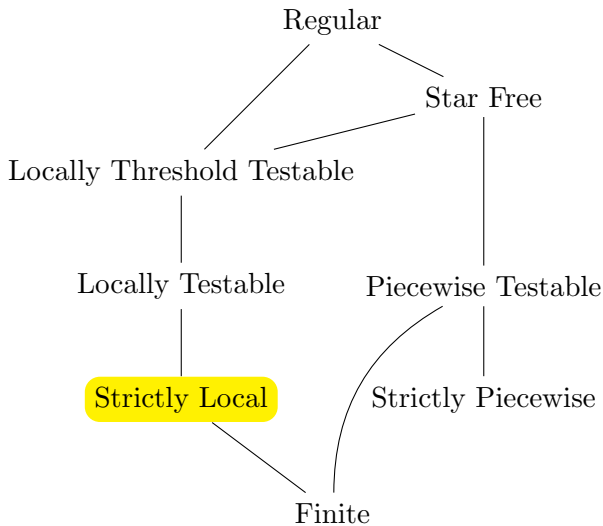
- To show  $L$  is not  $SL$ , find  $u_1, v_1, u_2, v_2$  and  $x$  for each  $k$ !

# EXERCISES

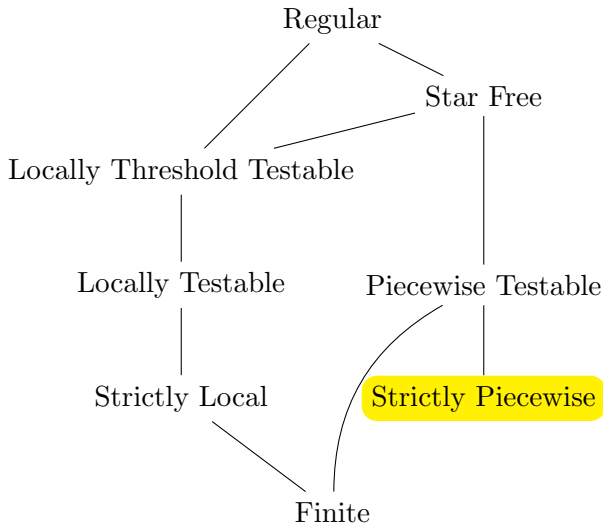
Prove the following languages are not SL for any  $k$ .

- (4) Strings contain at least one **b**.
- (5) Strings contain at most one **b**.
- (10) Strings contain an **a** between every pair of **bs**.

# SUBREGULAR COMPLEXITY



# SUBREGULAR COMPLEXITY



# SUBSEQUENCES

The notion of *subsequence* is key to understanding the “piecewise” subregular classes.

## k-Subsequences

- A string  $u$  is a subsequence of a string  $v = \sigma_1\sigma_2\ldots\sigma_n$  if

$$v \in \Sigma^* \sigma_1 \Sigma^* \sigma_2 \Sigma^* \ldots \Sigma^* \sigma_n \Sigma^* .$$

We write  $u \preceq v$ .

- $\text{subseq}_k(w)$  equals the set  $\{u \mid u \preceq_k w, |u| \leq k\}$ .
- $\text{subseq}_k(S)$  equals the set  $\bigcup_{w \in S} \preceq_k(w)$ .



# EXERCISE

Identify the following sets.

- 1  $\text{subseq}_2(abcd)$
- 2  $\text{subseq}_3(abcd)$
- 3  $\text{subseq}_2(aaaa, bbbb, caabc)$

# STRICTLY PIECEWISE LANGUAGES

A language  $L$  is Strictly  $k$ -Piecewise if there is a finite set of strings  $S \subseteq \mathbf{subseq}_k(\Sigma^*)$  such that

$$L = \{w \in \Sigma^* \mid \mathbf{subseq}_k(w) \subseteq S\}$$

A language  $L$  is Strictly Piecewise if there is some  $k$  for which it is Strictly  $k$ -Piecewise. We use  $\llbracket S \rrbracket$  to denote the extension of  $S$ .

# STRICTLY PIECEWISE LANGUAGES

## Notes

- Elements in  $S$  are called **permissible subsequences**.
- Elements in  $\text{subseq}_k(\Sigma^*)$  not in  $S$  are called **forbidden subsequences**.

## Exercises

- 1 Let  $S = \{\varepsilon, a, b, ab, ba\}$ . What is  $\llbracket S \rrbracket$ ?
- 2 For what  $k$  is  $\llbracket S \rrbracket$  Strictly  $k$ -Piecewise?
- 3 Which subsequences are forbidden?

# MORE EXERCISES

- 1 Recall this language mentioned earlier.

(5) Strings contain at most one **b**.

Explain why it is Strictly 2-Piecewise. Assume  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .

- 2 Recall Samala.

possible Samala words	impossible Samala words
ftojonowonowaf	stojonowonowaf
stojonowonowas	ftojonowonowas
pistonoskiwat	pisotonofikiwat
sanisotonoskiwas	faniapisotonofikiwas

Assuming  $\Sigma = \{\mathbf{s}, \mathbf{S}, \mathbf{t}, \mathbf{o}\}$ , what are the forbidden subsequences?

# STRICTLY PIECEWISE LANGUAGES

Intuitively, if  $L$  is Strictly  $k$ -Piecewise, then deciding whether a string  $w$  belongs to  $L$  simply requires checking its  $k$ -subsequences for forbidden ones. If any is found  $w$  is rejected. If every subsequence in  $w$  is permissible then  $w$  is accepted.

(Rogers et al. 2010)

# INFINITE HIERARCHY OF SP CLASSES

## Theorems

- ①  $\text{SP}_1 \subsetneq \text{SP}_2 \dots \text{SP}_k \subsetneq \text{SP}_{k+1} \dots \subseteq \text{SP}$
- ②  $\text{FIN} \not\subseteq \text{SP}$
- ③ For each  $k$ ,  $\text{SP}_k$  is closed under intersection, but neither complement nor union.

(Rogers et al. 2010)

# CHARACTERIZING STRICTLY PIECEWISE LANGUAGES

## Subsequence Closure

A language  $L$  is Strictly Piecewise iff whenever  $w \in L$  every subsequence of  $w$  also belongs to  $L$ .

Subsequence Closure helps us prove languages are not Strictly Piecewise.

(Rogers et al. 2010)

# EXERCISES

Prove the following languages are not SP.

- (4) Strings contain at least one **b**.
- (10) Strings contain an **a** between every pair of **bs**.
- (11) Strings contain an odd number of **bs**.



# SUBREGULAR COMPLEXITY

